

## Exam 1 from a Past Semester

1. Provide a brief answer to each of the following questions.

a) What do *perfect match* and *mismatch* mean in the context of Affymetrix GeneChip technology? Be as specific as possible in your answer.

b) True or False: Robots are used to print microarray slides so that gene locations can be easily randomized from slide to slide.

c) What is a probe set? Please be as specific as you can.

d) With a color depth of 16 bits/pixel, how many different signal values are possible?

e) Explain the meaning of *intensity-dependent dye bias*.

2. An experiment was conducted to study the effects of soil temperature on gene expression in developing soybean plants. A total of 18 soybean plants were randomly assigned to 6 containers so that each container held 3 plants. Each container had a separate control that could be used to adjust the soil temperature to any desired level. Three of the 6 containers were randomly selected to be set at a common cold soil temperature. The other 3 containers were set at a normal soil temperature. At the conclusion of the experiment, Affymetrix GeneChips were used to measure RNA levels with one GeneChip per plant.

- a) Name the experimental units in this experiment.
  
- b) Name the observational units in this experiment.
  
- c) Was blocking used in this experiment? If so, describe the blocks.
  
- d) Name the treatment factor(s) in this experiment and list the levels of each treatment factor.
  
- e) Write down a model for the data from a single gene. You may use the abbreviated notation described in class and in our course notes. Circle any terms in the model that you would treat as random.
  
- f) Suppose normalized natural-log-scale data for a single gene is as follows:

Container	Soil Temp.	Plant		
		1	2	3
1	normal	7	3	8
2	normal	6	3	3
3	normal	7	9	8
4	cold	1	4	1
5	cold	3	2	1
6	cold	5	7	3

(Note this “data” has been set at integer values to make computations easier.)

Provide an estimate of fold change for this gene that describes the effect of soil temperature on this gene’s expression level.

g) Provide a 95% confidence interval for the fold change estimated in part (f). Assume that the  $t$ -distribution quantile required for the computation of the confidence interval is 2.776.

h) Based on your 95% confidence interval, do you think the expression level of this gene was affected by soil temperature? Explain.

3. An experiment was conducted to study the effects of soil temperature and chemical treatment on gene expression in developing soybean plants. A total of 18 soybean plants were randomly assigned to 6 containers so that each container held 3 plants. Each container had a separate control that could be used to adjust the soil temperature to any desired level. Three of the 6 containers were randomly selected to be set at a cold soil temperature. The other 3 containers were set at a normal soil temperature. Three chemical treatments (A, B, and C) were randomly assigned to the plants in each container such that one of the three plants was selected for treatment with chemical A, another for treatment with chemical B, and the third for treatment with chemical C. At the conclusion of the experiment, Affymetrix GeneChips were used to measure RNA levels with one GeneChip per plant.

a) Name the experimental units in this experiment.

b) Write down a model for the data from a single gene. You may use the abbreviated notation described in class and in our course notes. Circle any terms in the model that you would treat as random.

c) Suppose that 18 two-color microarray slides will be used to measure expression instead of 18 GeneChips. Sketch a plot using our microarray circle and arrow notation to indicate how you would use two-color microarrays to measure RNA levels in the plants if the researchers are primarily interested in detecting differences among chemical treatments within each temperature.

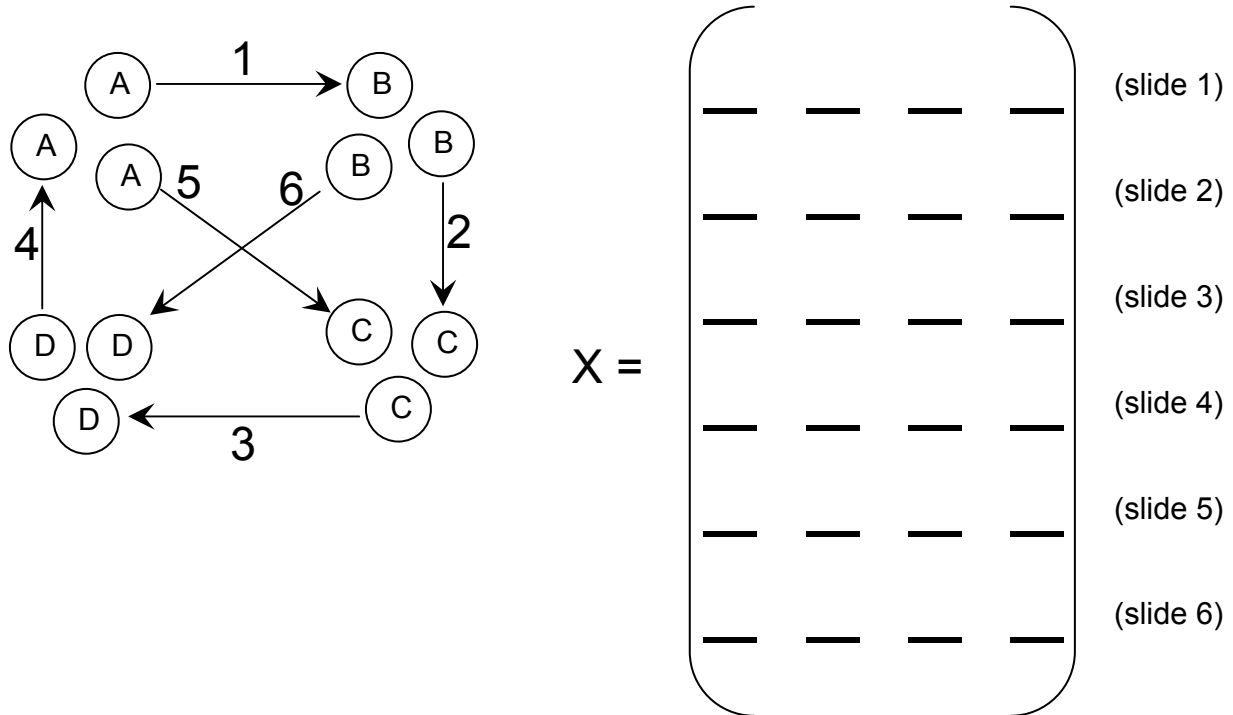
d) Write down a model for the two-color array data from a single gene based on the design you have specified above. You may use the abbreviated notation described in class and in our course notes. Circle any terms in the model that you would treat as random.

e) In the context of this example, explain the meaning of temperature-by-chemical interaction.

4. Suppose a two-color microarray experiment is to be conducted to compare the effect of four treatments (A, B, C, and D) on gene expression in maize. Suppose that treatment D is a control and that researchers are primarily interested in understanding which of the treatments A, B, and C differ from the control treatment D in terms of mean expression for each gene. The researchers have 12 experimental units and 6 slides available for the experiment. For any given gene, denote the mean expression of an observation as a function of dye and treatment according to the following table:

Treatment	Dye	Mean Expression
A	Cy3	$u+d_3+a$
B	Cy3	$u+d_3+b$
C	Cy3	$u+d_3+c$
D	Cy3	$u+d_3$
A	Cy5	$u+d_5+a$
B	Cy5	$u+d_5+b$
C	Cy5	$u+d_5+c$
D	Cy5	$u+d_5$

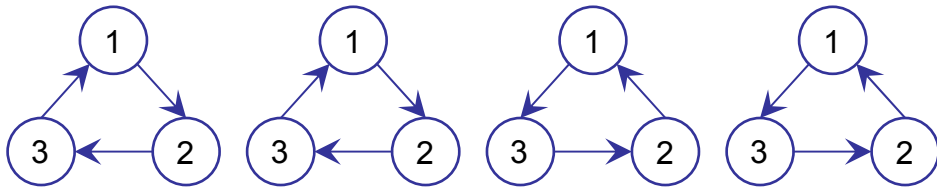
a) Consider the balanced incomplete block design (depicted below) that compares each treatment to each other treatment on exactly one slide. Provide the appropriate X matrix for this design. Assume that we will use the Cy5-Cy3 difference on each slide as our response variable and that our parameter vector for this analysis is  $[d_5-d_3, a, b, c]'$ . (Note that the numbers in the diagram below correspond to slide numbers. Please enter the rows of X accordingly.)



4. (continued)

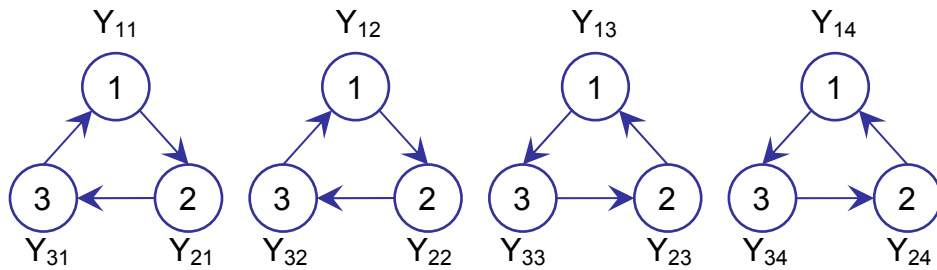
b) There exists at least one 6-slide 12-experimental-unit design that will dominate the design in part (a) for testing for differences between the control treatment (D) and each of the other three treatments (A, B, and C). Draw a diagram for a design that you believe will dominate the design in part (a). You will receive full credit if your chosen design does indeed dominate the design in (a). You do not need to do any calculation to show that your design dominates the design in (a); simply provide a diagram for such a design.

5. Consider a completely randomized experiment with three treatments denoted 1, 2, and 3 and four experimental units per treatment. Suppose mRNA levels are measured with two-color microarrays using the following design.



In class we discussed a mixed linear model for the 24 observations obtained for a single gene. This model included an overall mean  $\mu$  and fixed factors treatment and dye. Denote the treatment effects by  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ , and denote the dye effects by  $\delta_1$  and  $\delta_2$ . The model we discussed included variance components for the random factors slide, experimental unit, and residual. Denote these variance components by  $\sigma_s^2$ ,  $\sigma_u^2$ , and  $\sigma^2$  respectively. Assume this model is correct throughout this problem.

Suppose that instead of analyzing the 24 observations or instead of taking red-green differences as discussed in class, the two observations for each experimental unit are averaged to obtain a total of 12 averages denoted by  $Y_{ij}$  in the figure below.



a) Determine the expected value (mean) of  $Y_{11}$  in terms of the model parameters.

b) Determine the expected value (mean) of  $Y_{11} - Y_{21}$  in terms of model parameters.

c) Determine the variance of  $Y_{11} - Y_{21}$  in terms of the model parameters.